BSTA 620: Probability Lecture Notes

Di Shu, PhD

Department of Biostatistics, Epidemiology and Informatics University of Pennsylvania and Center for Pediatric Clinical Effectiveness Children's Hospital of Philadelphia Di.Shu@pennmedicine.upenn.edu

8th December 2021





Recap by case study

	drug A	drug B	
D+	a	b	m_1
D-	c	d	$-m_2$
	n_1	n_2	N

- Let $p_1 = a/n_1$ be the estimator for $\pi_1 = P(D + |drug = A)$, and $p_2 = b/n_2$ be the estimator for $\pi_2 = P(D + |drug = B)$
- Odds ratio $OR = \frac{\pi_1/(1 \pi_1)}{\pi_2/(1 \pi_2)}$, which can be estimated by $\widehat{OR} = \frac{p_1/(1 p_1)}{p_2/(1 p_2)}$

Shu, D (Penn&CHOP)

Recap by case study

- Why this estimator makes sense
- 95% confidence interval for \widehat{OR}
- Use and misuse of OR
- OR vs. risk ratio (RR)
- Another look at the rare disease condition by examining correlation



Definition 1.1.1: The set, *S*, of all possible outcomes of a particular experiment is called the **sample space** for the experiment

Definition 1.1.2: An **event** is any collection of possible outcomes of an experiment, that is, any subset of S (including S itself)

Definition 1.2.4: Given a sample space S and an associated sigma algebra \mathcal{B} , a **probability function** is a function P with domain \mathcal{B} that satisfies

- $P(A) \ge 0$ for all $A \in \mathcal{B}$
- P(S) = 1

If A₁, A₂,... ∈ B are pairwise disjoint, then P(∪_{i=1}[∞]A_i) = ∑_{i=1}[∞] P(A_i) (countable additivity)

• Items 1-3 are called Kolmogorov axioms of probability

The calculus of probabilities (consequences of 1.2.4)

Theorem 1.2.8: IF P is a probability function and A is any set in B, then

- $P(\emptyset) = 0$ where \emptyset is the empty set
- $P(A) \leq 1$

•
$$P(A^c) = 1 - P(A)$$



Theorem 1.2.9: If P is a probability function and A and B are any sets in \mathcal{B} , then

•
$$P(B \cap A^c) = P(B) - P(A \cap B)$$

•
$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

• If $A \subset B$, then $P(A) \leq P(B)$

$$P(\bigcup_{i=1}^{K} A_i) \le \sum_{i=1}^{K} P(A_i)$$

• Can see this from Theorem 1.2.9 which implies

 $P(A \cup B) \le P(A) + P(B)$

• Bonferroni's inequality tells us that if we make the probability of a type I error on any given comparison α/K , then the FWE will be

$$P(\bigcup_{i=1}^{K} \{ \text{type I error on test } i \}) \leq \sum_{i=1}^{K} \alpha/K = \alpha$$



Number of possible arrangements of size r from n objects

	Without replacement	With replacement
Ordered	$\frac{n!}{(n-r)!}$	n^r
Unordered	$\binom{n}{r}$	$\binom{n+r-1}{r}$



Definition 1.3.2: If A and B are events in S, and P(B) > 0, then the conditional probability of A given B, written as P(A|B), is

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Theorem 1.3.5 (Bayes's Rule): Let A_1, A_2, \ldots be a partition of the sample space, and let B be any set. Then for each $i = 1, 2, \ldots$,

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_{j=1}^{\infty} P(B|A_j)P(A_j)}$$

Definition 1.3.7: Two events A and B are statistically independent if

 $P(A \cap B) = P(A)P(B)$

Definition 1.3.12: A collection of events A_1, \ldots, A_n are **mutually independent** if for any subcollection A_{i_1}, \ldots, A_{i_k} , we have

$$P(\cap_{j=1}^{k} A_{i_j}) = \prod_{j=1}^{k} P(A_{i_j})$$



Theorem 1.3.9: If A and B are independent events, then the following pairs are also independent:

- $\bullet \ A \text{ and } B^c$
- $\bullet \ A^c \ {\rm and} \ B$
- $\bullet \ A^c \ {\rm and} \ B^c$
- Related If two random variables are independent then functions of those random variables are independent



Definition 1.4.1: A random variable is a function from a sample space S into the real numbers

• This is a simplified definition



Definition 1.5.1: the **cumulative distribution function or** *cdf* of a random variable X, denoted by $F_X(x)$, is defined by

 $F_X(x) = P_X(X \le x)$ for all x



Theorem 1.5.10: The following two statements are equivalent:

- $\bullet\,$ The random variables X and Y are identically distributed
- $F_X(x) = F_Y(x)$ for every x

Definition 1.6.1: The **probability mass function (pmf)** of a discrete random variable X is given by

$$f_X(x) = P(X = x)$$
 for all x

Definition 1.6.3: The probability density function (pdf) $f_X(x)$ of a continuous random variable X is the non-negative function that satisfies

$$F_X(x) = \int_{-\infty}^x f_X(t) dt$$
 for all x



Theorem 1.6.5: A function $f_X(x)$ is a pdf (or pmf) of a random variable X if and only if

• $f_X(x) \ge 0$ for all x

•
$$\sum_{x} f_X(x) = 1$$
 (pmf) or $\int_{-\infty}^{\infty} f_X(x) dx = 1$ (pdf)

- Common discrete distribution
 - Bernoulli
 - Binomial
 - Poisson (and Poisson Process)
 - Discrete uniform
 - Geometric
 - Hypergeometric
 - Negative binomial
 - Multinomial

• Common continuous distribution

- Normal
- Cauchy
- Uniform
- Exponential
- Gamma
- Weibull
- Beta
- log normal
- Double exponential
- Chi-squared, t, F
- And their relations

Theorem 2.1.5: Let X have pdf $f_X(x)$ and let Y = g(X), where g is a monotone function. Let \mathcal{X} and \mathcal{Y} be defined as in Theorem 2.1.3. Suppose that $f_X(x)$ is continuous on \mathcal{X} and that $g^{-1}(y)$ has a continuous derivative on \mathcal{Y} . Then the pdf of Y is given by

$$f_Y(y) = \begin{cases} f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right| & y \in \mathcal{Y} \\ 0 & otherwise \end{cases}$$

Theorem 2.1.8: Let X have pdf $f_X(x)$ and let Y = g(X), and define \mathcal{X} as above. Suppose there exists a partition A_0, A_1, \ldots, A_k , of \mathcal{X} such that $P(X \in A_0) = 0$ and $f_X(x)$ is continuous on each A_i . Further, suppose there exist functions $g_1(x), \ldots, g_k(x)$, defined on A_1, \ldots, A_k , respectively, satisfying

- $g(x) = g_i(x)$ for $x \in A_i$
- $g_i(x)$ is monotone on A_i
- The set $\mathcal{Y} = \{y : y = g_i(x) \text{ for some } x \in A_i\}$ is the same for each $i = 1, \dots, k$, and
- $g_i^{-1}(y)$ has a continuous derivative on \mathcal{Y} , for each $i = 1, 2, \ldots, k$. Then

$$f_{Y}(y) = \begin{cases} \sum_{i=1}^{k} f_{X}(g_{i}^{-1}(y)) \left| \frac{d}{dy} g_{i}^{-1}(y) \right| & y \in \mathcal{Y} \\ 0 & otherwise \end{cases}$$

Shu, D (Penn&CHOP)



Theorem 2.1.10 (Probability integral transformation): Let X have continuous cdf $F_X(x)$ and define the random variable Y as $Y = F_X(X)$. Then $Y \sim uniform(0, 1)$, that is, $P(Y \le y) = y$, 0 < y < 1



Definition 2.2.1: The **expected value or mean** of a random variable g(X), denoted $E\{g(X)\}$, is

$$E\{g(X)\} = \begin{cases} \int_{-\infty}^{\infty} g(x) f_X(x) dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} g(x) f_X(x) = \sum_{x \in \mathcal{X}} g(x) P(X = x) & \text{if } X \text{ is discrete} \end{cases}$$

provided that the integral or sum exists. If $E|g(X)| = \infty$, we say that $E\{g(X)\}$ does not exist

• Consider g(X) = X:

Definition: The **expected value or mean** of a random variable X, denoted E(X), is

 $E(X) = \begin{cases} \int_{-\infty}^{\infty} x f_X(x) dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} x f_X(x) = \sum_{x \in \mathcal{X}} x P(X = x) & \text{if } X \text{ is discrete} \end{cases}$

provided that the integral or sum exists. If ${\rm E}|X|=\infty,$ we say that ${\rm E}(X)$ does not exist

Definition 2.3.1: For each integer n, the nth **moment** of X (or $F_X(x)$), μ'_n , is

 $\mu'_n = \mathcal{E}(X^n)$

The *n*th central moment of *X*, μ_n , is

$$\mu_n = \mathrm{E}\{(X - \mu)^n\}$$

where $\mu = \mu'_1 = \mathcal{E}(X)$

Definition 2.3.2: The **variance** of a random variable X is its second central moment, $Var(X) = E[{X - E(X)}^2]$. The positive square root of Var(X) is the **standard deviation** of X

• $Var(X) = E(X^2) - {E(X)}^2$

Theorem 2.2.5: Let X be a random variable and let a, b and c be constants. Then for any functions $g_1(x)$ and $g_2(x)$ whose expectations exist,

- $E(ag_1(X) + bg_2(X) + c) = aE\{g_1(X)\} + bE\{g_2(X)\} + c$
- If $g_1(x) \ge 0$ for all x, then $\mathrm{E}\{g_1(X)\} \ge 0$
- If $g_1(x) \ge g_2(x)$ for all x, then $\mathrm{E}\{g_1(X)\} \ge \mathrm{E}\{g_2(X)\}$
- If $a \leq g_1(x) \leq b$ for all x, then $a \leq \mathrm{E}\{g_1(X)\} \leq b$

Theorem 2.3.4: If X is a random variable with finite variance, then for any constants a and b,

$$\operatorname{Var}(aX+b) = a^2 \operatorname{Var}(X)$$



Definition 2.3.6: Let X be a random variable with cdf F_X . The **moment** generating function (mgf) of X (or of F_X), denoted by $M_X(t)$, is

$$M_X(t) = \mathcal{E}(e^{tX})$$

provided that the expectation exists for t in some neighborhood of 0 (i.e. $\exists h > 0$ such that, $\forall t \in (-h, h)$, $E(e^{tX})$ exists). If the expectation does not exist in a neighborhood of 0, we say that the mgf does not exist

Theorem 2.3.7: If X has mgf $M_X(t)$, then

 $\mathcal{E}(X^n) = M_X^{(n)}(0)$

where we define

$$M_X^{(n)}(0) = \frac{d^n}{dt^n} M_X(t)\Big|_{t=0}$$

That is, the *n*th moment is equal to the *n*th derivative of $M_X(t)$ evaluated at t = 0

Theorem 2.3.11: Let $F_X(x)$ and $F_Y(y)$ be two cdfs all of whose moments exist

- If X and Y have bounded support, then $F_X(u) = F_Y(u)$ for all u if and only if $E(X^r) = E(Y^r)$ for all integers r = 0, 1, 2, ...
- If the mgfs exist and $M_X(t) = M_Y(t)$ for all t in some neighborhood of 0, then $F_X(u) = F_Y(u)$ for all u

Theorem 2.3.12 (Convergence of mgfs): Suppose $\{X_i, i = 1, 2, ...\}$ is a sequence of random variables, each with mgf $M_{X_i}(t)$. Furthermore, suppose that

 $\lim_{i \to \infty} M_{X_i}(t) = M_X(t)$

for all t in a neighborhood of 0, and $M_X(t)$ is an mgf. Then there is a unique cdf F_X whose moments are determined by $M_X(t)$ and, for all x where $F_X(x)$ is continuous, we have

$$\lim_{\to\infty} F_{X_i}(x) = F_X(x)$$

• That is, convergence of mgfs to an mgf (for |t| < h) implies convergence of the cdfs

Shu, D (Penn&CHOP)

Theorem 2.3.15: For any constants a and b, the mgf of the random variable aX + b is given by

 $M_{aX+b}(t) = e^{bt} M_X(at)$

• **Exponential families** A family of pdfs or pmfs is called an exponential family if it can be expressed as

$$f(x; \theta) = h(x)c(\theta) \exp\left(\sum_{i=1}^{k} w_i(\theta)t_i(x)\right)$$

where $h(x) \ge 0$ and $t_1(x), \ldots, t_k(x)$ are real-valued functions of the observation x (they cannot depend on θ), and $c(\theta) \ge 0$ and $w_1(\theta), \ldots, w_k(\theta)$ are real-valued functions of the possibly vector-valued parameter θ (they cannot depend on x)

Definition 3.5.5: Let f(x) be any pdf. Then for any $\mu \in (-\infty, \infty)$ and any $\sigma > 0$, the family of pdfs $(1/\sigma)f((x-\mu)/\sigma)$, indexed by the parameter (μ, σ) , is called the **location-scale family with standard pdf** f(x); μ is called the **location parameter** and σ is called the **scale parameter**



Definition 4.1.1: An *n*-dimensional random vector is a function from a sample space S into \mathcal{R}^n , *n*-dimensional Euclidean space
Definition 4.1.3: Let (X, Y) be a discrete bivariate random vector. Then function f(x, y) from \mathcal{R}^2 into \mathcal{R} defined by f(x, y) = P(X = x, Y = y) is called the **joint probability mass function or joint pmf** of (X, Y)

- To be clearer, the notation $f_{X,Y}(x,y)$ will be used
- Use the joint pmf to calculate probability of any event:

$$P\{(X,Y)\in A\} = \sum_{(x,y)\in A} f(x,y)$$

Theorem 4.1.6: Let (X, Y) be a discrete bivariate random vector with joint pmf $f_{X,Y}(x, y)$. Then the marginal pmfs of X and Y, $f_X(x)$ and $f_Y(y)$, respectively, are given by

$$f_X(x) = \sum_{y \in \mathcal{R}} f_{X,Y}(x,y)$$

and

$$f_Y(y) = \sum_{x \in \mathcal{R}} f_{X,Y}(x,y)$$

Definition 4.1.10: A function f(x, y) from \mathcal{R}^2 into \mathcal{R} is called a **joint probability density function or joint pdf** of the continuous bivariate random vector (X, Y) if, for every $A \subset \mathcal{R}^2$

$$P((X,Y) \in A) = \int_A \int f(x,y) dx dy$$

• Expectation of a function of a bivariate continuous random variable

$$\mathbb{E}\{g(X,Y)\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x,y)f(x,y)dxdy$$

• The marginal probability density functions of X and Y are given by

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy, \ -\infty < x < \infty$$

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx, \ -\infty < y < \infty$$

• Any function $f(x,y) \geq 0$, $\forall (x,y) \in \mathcal{R}^2$ with

$$1 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy$$

is the joint pdf of some continuous bivariate random vector (X, Y)

Definition 4.2.1: Let (X, Y) be a discrete bivariate random vector with joint pmf f(x, y) and marginal pmfs $f_X(x)$ and $f_Y(y)$. For any x such that $P(X = x) = f_X(x) > 0$, the **conditional pmf of** Y **given that** X = x is the function of y denoted by f(y|x) and defined by

$$f(y|x) = P(Y = y|X = x) = \frac{f(x,y)}{f_X(x)}$$

For any y such that $P(Y = y) = f_Y(y) > 0$, the conditional pmf of X given that Y = y is the function of x denoted by f(x|y) and defined by

$$f(x|y) = P(X = x|Y = y) = \frac{f(x,y)}{f_Y(y)}$$

• Define conditional expected value of g(Y) given X = x as

 $\mathbf{E}\{g(Y)|x\} = \sum_{y} g(y)f(y|x)$

Definition 4.2.3: Let (X, Y) be a continuous bivariate random vector with joint pdf f(x, y) and marginal pdfs $f_X(x)$ and $f_Y(y)$. For any x such that $f_X(x) > 0$, the **conditional pdf of** Y **given that** X = x is the function of y denoted by f(y|x) and defined by

$$f(y|x) = \frac{f(x,y)}{f_X(x)}$$

For any y such that $f_Y(y) > 0$, the **conditional pdf of** X **given that** Y = y is the function of x denoted by f(x|y) and defined by

$$f(x|y) = \frac{f(x,y)}{f_Y(y)}$$

• Define conditional expected value of g(Y) given X = x as

$$\mathbf{E}\{g(Y)|x\} = \int_{-\infty}^{\infty} g(y)f(y|x)dy$$

• The variance of the probability distribution described by f(y|x), denoted by Var(Y|x), is called the **conditional variance** of Y given X = x

$$Var(Y|x) = E(Y^2|x) - {E(Y|x)}^2$$

Definition 4.2.5: Let (X, Y) be a bivariate random vector with joint pdf or pmf f(x, y) and marginal pdfs or pmfs $f_X(x)$ and $f_Y(y)$. Then X and Y are called **independent random variables** if for every $x \in \mathcal{R}$ and $y \in \mathcal{R}$,

 $f(x,y) = f_X(x)f_Y(y)$

Lemma 4.2.7: Let (X, Y) be a bivariate random vector with joint pdf or pmf f(x, y). Then X and Y are **independent random variables** if and only if there exist functions g(x) and h(y) such that, for every $x \in \mathcal{R}$ and $y \in \mathcal{R}$,

f(x,y) = g(x)h(y)

Theorem 4.2.10: Let X and Y be independent random variables.

• For any
$$A \subset \mathcal{R}$$
 and $B \subset \mathcal{R}$,

$$P(X \in A, Y \in B) = P(X \in A)P(Y \in B)$$

that is, the events $\{X \in A\}$ and $\{Y \in B\}$ are **independent events** • Let g(x) be a function only of x and h(y) be a function only of y. Then

$$\mathbf{E}\{g(X)h(Y)\} = \mathbf{E}\{g(X)\} \cdot \mathbf{E}\{h(Y)\}$$

Theorem 4.2.12: Let X and Y be independent random variables with moment generating functions $M_X(t)$ and $M_Y(t)$. Then the moment generating function of the random variable Z = X + Y is given by

 $M_Z(t) = M_X(t)M_Y(t)$

Theorem 4.2.14: Let $X \sim n(\mu, \sigma^2)$ and $Y \sim n(\gamma, \tau^2)$ be independent normal random variables. Then the random variable Z = X + Y has a $n(\mu + \gamma, \sigma^2 + \tau^2)$ distribution

- Bivariate transformations (see previous notes for details)
 - Discrete case

$$f_{U,V}(u,v) = P(U = u, V = v) = P((X,Y) \in A_{uv}) = \sum_{(x,y) \in A_{uv}} f_{X,Y}(x,y)$$

• Continous case (assuming 1 to 1 transformation)

$$f_{U,V}(u,v) = \begin{cases} f_{X,Y}(h_1(u,v),h_2(u,v))|J| & (u,v) \in \mathcal{B} \\ 0 & otherwise \end{cases}$$

• Continous case (not 1 to 1)

$$f_{U,V}(u,v) = \sum_{i=1}^{k} f_{X,Y}(h_{1i}(u,v), h_{2i}(u,v))|J_i|$$

Theorem 4.3.5: Let X and Y be independent random variables. Let g(X) be a function only of X and h(Y) be a function only of Y. Then the random variables U = g(X) and V = h(Y) are independent

Theorem 4.4.3: If X and Y are any two random variables, then

 $\mathcal{E}(X) = \mathcal{E}\{\mathcal{E}(X|Y)\}$

provided that the expectations exist

Theorem 4.4.7 (Conditional variance identity): For any two random variables X and Y,

$$\operatorname{Var}(X) = \operatorname{E}\{\operatorname{Var}(X|Y)\} + \operatorname{Var}\{\operatorname{E}(X|Y)\}$$

provided that the expectations exist

Shu, D (Penn&CHOP)

Definition 4.5.1: The **covariance** of X and Y is the number defined by

$$Cov(X, Y) = E\{(X - \mu_X)(Y - \mu_Y)\}\$$

Theorem 4.5.3: $Cov(X, Y) = E(XY) - \mu_X \mu_Y$

Definition 4.5.2: The **correlation** of *X* and *Y* is the number defined by

$$\rho_{X,Y} = \frac{\operatorname{Cov}(X,Y)}{\sigma_X \sigma_Y}$$

The value $\rho_{X,Y}$ is also called the **correlation coeffcient**

Shu, D (Penn&CHOP)

54 / 90

Theorem 4.5.5: If X and Y are independent random variables, then Cov(X, Y) = 0 and $\rho_{X,Y} = 0$

• Independence implies 0 covariance, but not the other way around



Theorem 4.5.6: If X and Y are any two random variables and a and b are any two constants, then

$$\operatorname{Var}(aX + bY) = a^{2}\operatorname{Var}(X) + b^{2}\operatorname{Var}(Y) + 2ab\operatorname{Cov}(X, Y)$$

If X and Y are independent random variables, then

$$Var(aX + bY) = a^{2}Var(X) + b^{2}Var(Y)$$

Theorem 4.5.7: For any random variables X and Y,

•
$$-1 \le \rho_{XY} \le 1$$

• $|\rho_{XY}| = 1$ if and only if there exist numbers $a \neq 0$ and b such that P(Y = aX + b) = 1. If $\rho_{XY} = 1$, then a > 0, and if $\rho_{XY} = -1$, then a < 0

• Multivariate distributions (see previous notes for details)

- joint distribution
- marginal distribution
- conditional distribution
- expected value

Definition 4.6.5: Let X_1, \ldots, X_n be random vectors with joint pdf or pmf $f(x_1, \ldots, x_n)$. Let $f_{X_i}(x_i)$ denote the marginal pdf or pmf of X_i . Then X_1, \ldots, X_n are called **mutually independent random vectors** if, for every (x_1, \ldots, x_n) ,

$$f({m x}_1,\ldots,{m x}_n)=f_{{m X}_1}({m x}_1)\cdots f_{{m X}_n}({m x}_n)=\prod_{i=1}^n f_{{m X}_i}({m x}_i)$$

• If X_i s are all one-dimensional, then X_1, \ldots, X_n are called **mutually** independent random variables

Theorem 4.6.6 (Generalization of Theorem 4.2.10): Let X_1, \ldots, X_n be **mutually independent random variables**. Let g_1, \ldots, g_n be real-valued functions such that $g_i(x_i)$ is a function only of x_i , $i = 1, \ldots, n$. Then

 $\mathbf{E}\{g_1(X_1)\cdots g_n(X_n)\}=\mathbf{E}\{g_1(X_1)\}\cdots \mathbf{E}\{g_n(X_n)\}$

Theorem 4.6.7 (Generalization of Theorem 4.2.12): Let X_1, \ldots, X_n be **mutually independent random variables** with mgfs $M_{X_1}(t), \ldots, M_{X_n}(t)$. Let $Z = X_1 + \cdots + X_n$. Then the mgf of Z is

$$M_Z(t) = M_{X_1}(t) \cdots M_{X_n}(t)$$

In particular, if X_1, \ldots, X_n all have the same distribution with mgf $M_X(t)$, then

 $M_Z(t) = \{M_X(t)\}^n$

Theorem 4.6.11 (Generalization of Lemma 4.2.7): Let X_1, \ldots, X_n be random vectors. Then X_1, \ldots, X_n are mutually independent random vectors if and only if there exist functions $g_i(x_i)$, $i = 1, \ldots, n$, such that the joint pdf or pmf of (X_1, \ldots, X_n) can be written as

$$f(oldsymbol{x}_1,\ldots,oldsymbol{x}_n)=g_1(oldsymbol{x}_1)\cdots g_n(oldsymbol{x}_n)$$

Theorem 4.6.12 (Generalization of Theorem 4.3.5): Let X_1, \ldots, X_n be independent random vectors. Let $g_i(x_i)$ be a function only of x_i , $i = 1, \ldots, n$. Then the random variables $U_i = g_i(X_i)$, $i = 1, \ldots, n$, are mutually independent

- Chebyshev's inequality
- Markov inequality
- Normal tail probability
- Hölder's inequality
- Cauchy-Schwarz inequality
- Covariance inequality
- Minkowski's inequality
- Jensen's inequality
- Inequality for means

Definition 5.1.1: The random variables X_1, \ldots, X_n are called a **random** sample of size n from the population f(x) if X_1, \ldots, X_n are mutually independent random variables and the marginal pdf or pmf of each X_i is the same function f(x). Alternatively, X_1, \ldots, X_n are called **independent and identically distributed random variables** with pdf or pmf f(x). This is commonly abbreviated to **iid** random variables

Definition 5.2.1: Let X_1, \ldots, X_n be a random sample of size n from a population and let $T(x_1, \ldots, x_n)$ be a real-valued or vector-valued function whose domain includes the sample space of (X_1, \ldots, X_n) . Then the random variable or random vector $Y = T(X_1, \ldots, X_n)$ is called a **statistic**. The probability distribution of a statistic Y is called the **sampling distribution of** Y

Definition 5.2.2: The **sample mean** is the arithmetic average of the values in a random sample. It is usually denoted by

$$\bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$$

Definition 5.2.3: The sample variance is the statistic defined by

$$S^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (X_{i} - \bar{X})^{2}$$

The sample standard deviation is the statistic defined by $S = \sqrt{S^2}$

Shu, D (Penn&CHOP)

Theorem 5.2.6: Let X_1, \ldots, X_n be a random sample from a population with mean μ and variance $\sigma^2 < \infty$. Then

- $E(\bar{X}) = \mu$ (unbiasedness)
- $\operatorname{Var}(\bar{X}) = \frac{\sigma^2}{n}$ $\operatorname{E}(S^2) = \sigma^2$ (unbiasedness)

- Important factoids Let X_1, \ldots, X_n be random variables whose expectations and variances exist
 - $\operatorname{E}(X_1 + \dots + X_n) = \sum_{i=1}^n \operatorname{E}(X_i)$
 - Var $(X_1 + \dots + X_n) = \sum_{i=1}^n \operatorname{Var}(X_i) + \sum_{i \neq j} \operatorname{Cov}(X_i, X_j)$
- Note for a random sample,

$$\operatorname{Var}(X_1 + \dots + X_n) = n\operatorname{Var}(X_1)$$

• Note for a previous example of sampling without replacement,

$$\operatorname{Var}(X_1 + \dots + X_n) = n\operatorname{Var}(X_1) + n(n-1)\operatorname{Cov}(X_1, X_2)$$

Theorem 5.2.7: Let X_1, \ldots, X_n be a random sample from a population with mgf $M_X(t)$. Then the mgf of the sample mean is

 $M_{\bar{X}}(t) = \{M_X(t/n)\}^n$

• When the mgf of \bar{X} is not recognizable, or the population mgf does not exist, the transformation method might be used. In such cases, the following **convolution formula** is useful

Theorem 5.2.9: If X and Y are **independent**, continuous random variables with pdfs $f_X(x)$ and $f_Y(y)$, then the pdf of Z = X + Y is

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(w) f_Y(z-w) dw$$

Theorem 5.3.1: Let X_1, \ldots, X_n be a random sample from a $n(\mu, \sigma^2)$ distribution, and let $\bar{X} = (1/n) \sum_{i=1}^n X_i$ and $S^2 = \{1/(n-1)\} \sum_{i=1}^n (X_i - \bar{X})^2$. Then

• \bar{X} and S^2 are independent random variables

•
$$\bar{X} \sim n(\mu, \sigma^2/n)$$

•
$$(n-1)S^2/\sigma^2 \sim \chi^2_{n-1}$$


Definition 5.4.1: The order statistics of a random sample X_1, \ldots, X_n are the sample values placed in ascending order. They are denoted by $X_{(1)}, \ldots, X_{(n)}$

Theorem 5.4.4: Let $X_{(1)}, \ldots, X_{(n)}$ denote the order statistics of a random sample, X_1, \ldots, X_n , from a continuous population with cdf $F_X(x)$ and pdf $f_X(x)$. Then the pdf of $X_{(j)}$ is

$$f_{X_{(j)}}(x) = \frac{n!}{(j-1)!(n-j)!} f_X(x) \{F_X(x)\}^{j-1} \{1 - F_X(x)\}^{n-j}$$

• Proof: $F_{X_{(j)}}(x) = \sum_{k=j}^{n} {n \choose k} \{F_X(x)\}^k \{1 - F_X(x)\}^{n-k}$ and then differentiate

Theorem 5.4.6: Let $X_{(1)}, \ldots, X_{(n)}$ denote the order statistics of a random sample, X_1, \ldots, X_n , from a continuous population with cdf $F_X(x)$ and pdf $f_X(x)$. Then the joint pdf of $X_{(i)}$ and $X_{(j)}$, $1 \le i < j \le n$, is

$$= \frac{f_{X_{(i)},X_{(j)}}(u,v)}{(i-1)!(j-1-i)!(n-j)!} f_X(u) f_X(v) \{F_X(u)\}^{i-1} \times \{F_X(v) - F_X(u)\}^{j-1-i} \{1 - F_X(v)\}^{n-j}$$

for $-\infty < u < v < \infty$

 $X_n \xrightarrow{p} X$: A sequence of random variables, X_1, X_2, \ldots , converges in **probability** to a random variable X if, for every $\epsilon > 0$,

 $\lim_{n \to \infty} P(|X_n - X| < \epsilon) = 1$

 $X_n \xrightarrow{a.s} X$: A sequence of random variables, X_1, X_2, \ldots , converges almost surely to a random variable X if, for every $\epsilon > 0$,

$$P(\lim_{n \to \infty} |X_n - X| < \epsilon) = 1$$



Definition 5.5.10: A sequence of random variables, X_1, X_2, \ldots converges in distribution to a random variable X if

$$\lim_{n \to \infty} F_{X_n}(x) = F_X(x)$$

at all points x where $F_X(x)$ is continuous. We also write this

$$X_n \xrightarrow{d} X$$



Theorem 5.5.2 (Weak Law of Large Numbers): Let X_1, X_2, \ldots be iid random variables with $E(X_i) = \mu$ and $Var(X_i) = \sigma^2 < \infty$. Define $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. Then, for every $\epsilon > 0$,

 $\lim_{n \to \infty} P(|\bar{X}_n - \mu| < \epsilon) = 1$

That is, \bar{X}_n converges in probability to μ :

 $\bar{X}_n \xrightarrow{p} \mu$

We say that \overline{X}_n is **consistent** for μ

Shu, D (Penn&CHOP)

Theorem 5.5.9 (Strong Law of Large Numbers): Let X_1, X_2, \ldots be iid random variables with $E(X_i) = \mu$ and $Var(X_i) = \sigma^2 < \infty$. Define $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. Then, for every $\epsilon > 0$,

$$P(\lim_{n \to \infty} |\bar{X}_n - \mu| < \epsilon) = 1$$

That is,

$$\bar{X}_n \xrightarrow{a.s.} \mu$$

Theorem 5.5.14 (Central Limit Theorem): Let X_1, X_2, \ldots be a sequence of iid random variables whose mgfs exist in a neighborhood of 0 (i.e. $M_{Xi}(t)$ exists for |t| < h, for some positive h). Let $E(X_i) = \mu$ and $Var(X_i) = \sigma^2 > 0$ (Both μ and σ^2 are finite because the mgf exists). Define $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. Let $G_n(x)$ denote the cdf of $\sqrt{n}(\bar{X}_n - \mu)/\sigma$. Then, for any $x \in (-\infty, \infty)$,

$$\lim_{n \to \infty} G_n(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy$$

That is,

$$\sqrt{n}(\bar{X}_n - \mu) / \sigma \xrightarrow{d} \mathbf{n}(0, 1)$$

Theorem 5.5.15 (Stronger form of the CLT): Let X_1, X_2, \ldots be a sequence of iid random variables with $E(X_i) = \mu$ and $0 < Var(X_i) = \sigma^2 < \infty$. Define $\bar{X}_n = (1/n) \sum_{i=1}^n X_i$. Let $G_n(x)$ denote the cdf of $\sqrt{n}(\bar{X}_n - \mu)/\sigma$. Then, for any $x \in (-\infty, \infty)$,

$$\lim_{n \to \infty} G_n(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy$$

That is,

$$\sqrt{n}(\bar{X}_n - \mu) / \sigma \stackrel{d}{\longrightarrow} \mathbf{n}(0, 1)$$

- Continuous mapping theorem/Mann-Wald mapping theorem
 Suppose that X₁, X₂,... is a sequence of random variables and f : R → R is a Borel function (includes continuous functions) whose set D of discontinuities is such that ω : X(ω) ∈ D ∈ F and P(X ∈ D) = 0. If X_n converges to X either
 - almost surely
 - in-probability, or
 - in distribution

Then $g(X_n)$ converges to g(X) in the same sense

Theorem 5.5.17 (Slutsky's Theorem): If $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{p} a$, where a is a constant, then

- $X_n Y_n \xrightarrow{d} aX$
- $X_n + Y_n \xrightarrow{d} X + a$

Theorem 5.5.24 (Delta Method): Let Y_n be a sequence of random variables that satisfies $\sqrt{n}(Y_n - \theta) \stackrel{d}{\longrightarrow} n(0, \sigma^2)$. For a given function g and a specific value of θ , suppose that $g'(\theta)$ exists and is not 0. Then

$$\sqrt{n}\{g(Y_n) - g(\theta)\} \stackrel{d}{\longrightarrow} \mathbf{n}(0, \sigma^2\{g'(\theta)\}^2)$$

Theorem 5.5.26 (Second-order Delta Method): Let Y_n be a sequence of random variables that satisfies $\sqrt{n}(Y_n - \theta) \stackrel{d}{\longrightarrow} n(0, \sigma^2)$. For a given function g and a specific value of θ , suppose that $g'(\theta) = 0$ and $g''(\theta)$ exists and is not 0. Then

$$n\{g(Y_n) - g(\theta)\} \xrightarrow{d} \sigma^2 \left\{ \frac{g''(\theta)}{2} \right\} \chi_1^2$$

• Let $T = (T_1, \ldots, T_k)^T$ with mean μ and variance-covariance matrix Σ_T . Let $Q = G(T) = (g_1(T), \ldots, g_m(T))^T$. Then

 $\mathrm{E}(\boldsymbol{Q}) \approx \boldsymbol{G}(\boldsymbol{\mu}) \quad \text{and} \quad \mathrm{Var}(\boldsymbol{Q}) \approx H(\boldsymbol{\mu}) \Sigma_{\boldsymbol{T}} H(\boldsymbol{\mu})^T$ where $H(\boldsymbol{\mu}) = H(\boldsymbol{t})|_{\boldsymbol{t}=\boldsymbol{\mu}}$ and

$$H(oldsymbol{t}) = egin{bmatrix} rac{\partial g_1(oldsymbol{t})}{\partial t_1} & \ldots & rac{\partial g_1(oldsymbol{t})}{\partial t_k} \ dots & \ddots & dots \ rac{\partial g_m(oldsymbol{t})}{\partial t_1} & \ldots & rac{\partial g_m(oldsymbol{t})}{\partial t_k} \end{bmatrix}$$

- Final exam focuses on materials after the midterm (after Section 3.3)
- A couple of notes (1/3)
 - $E(\cdot)$
 - is sum or integration includes mean, variance, moments and mgfs
 - handy calculation of $\mathrm{E}(X)$ and $\mathrm{Var}(X)$ using hierarchy
 - Derive marginal and conditional pdfs/pmfs from a joint pdf/pmf
 - Univariate, bivariate, or multivariate transformations
 - discrete case or continuous case?
 - do we have 1 to 1 transformation?

- A couple of notes (2/3)
 - mgfs can be used to
 - derive moments
 - derive variance through 1st and 2nd moments
 - name/recognize a distribution (note: simplified derivation for sum of independent r.v.s)
 - examine convergence in distribution
 - Independence of r.v.s
 - is verified when both joint density and support region factor
 - implies the independence of events and functions of r.v.s.
 - implies that, expectation of product = product of expectation
 - implies that covariance and correlation = 0; the opposite direction is true for normal r.v.s but does not generally hold

- A couple of notes (3/3)
 - Relations between r.v.s (e.g. square of standard normal is chi-square)
 - Random sample means iid
 - Sample mean and sample variance are unbiased. More properties when assuming normal distribution
 - Order statistics (min, max and more)
 - Probabilistic inequalities and applications
 - Three modes of convergence and their relations
 - How to apply LLN, CLT and Delta method and when?

- Always double check assumptions when applying any theorems or properties (e.g. $M_{X+Y}(t) = M_X(t)M_Y(t)$ require X and Y be independent)
- Familiarity with common distributions and their relations is useful and sometimes can save time
- Useful tools: proof by contradiction, induction, recursive relation, recognize a kernel, integration (dxdy or dydx, polar coordinates), geometric argument (e.g. find the area), inequalities, Taylor expansion, etc.
- Pay attention to support region
- Go back to definitions if no clue; they are fundamental and often the first step of solutions